



Lorenz-Spreen, P., Lewandowsky, S., Sunstein, C. R., & Hertwig, R. (2020). How behavioural sciences can promote truth, autonomy and democratic discourse online. *Nature Human Behaviour*, 2020. <https://doi.org/10.1038/s41562-020-0889-7>

Peer reviewed version

Link to published version (if available):  
[10.1038/s41562-020-0889-7](https://doi.org/10.1038/s41562-020-0889-7)

[Link to publication record in Explore Bristol Research](#)  
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Nature Research at <https://www.nature.com/articles/s41562-020-0889-7> . Please refer to any applicable terms of use of the publisher.

## University of Bristol - Explore Bristol Research

### General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available: <http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

## How behavioural sciences can promote truth, autonomy and democratic discourse online

Deleted: and

Philipp Lorenz-Spreen<sup>1\*</sup>, Stephan Lewandowsky<sup>2,3</sup>, Cass R. Sunstein<sup>4</sup>, Ralph Hertwig<sup>1</sup>

<sup>1</sup> Center for Adaptive Rationality, Max Planck Institute for Human Development, Berlin, Germany

<sup>2</sup> School of Psychological Science and Cabot Institute, University of Bristol, Bristol, United Kingdom

<sup>3</sup> School of Psychological Science, University of Western Australia, Perth, Western Australia, Australia

<sup>4</sup> Harvard Law School, Cambridge, Massachusetts, United States

\*Correspondence to: lorenz-spreen@mpib-berlin.mpg.de

**Abstract:** Public opinion is shaped by online content, spread via social media and curated

Deleted: in large part

algorithmically. The current online ecosystem has been designed predominantly to capture user attention rather than promote deliberate cognition and autonomous choice. The resultant

Deleted: by algorithmic recommender systems

information overload, finely tuned personalization, and distorted social cues, in turn, pave the

Deleted: ;

way for manipulation and the spread of false information. How can transparency and autonomy

be promoted instead, thus fostering the positive potential of the Web? Effective Web governance

Deleted: , to take advantage of

informed by behavioural research is critically needed to empower individuals online. We identify

Deleted: We argue that effective

technologically available yet largely untapped cues that can be harnessed to indicate the

Deleted: in the online world

epistemic quality of online content, the factors underlying algorithmic decisions and the degree

Deleted: behind

of consensus in online debates. We then map out two classes of behavioural interventions—

Deleted: social media

nudging and boosting—that enlist these cues to redesign online environments for informed and

Deleted: in ways that foster user autonomy

autonomous choice and democratic discourse.

Deleted: promote truth.

40 To the extent that a “wealth of information creates a poverty of attention” (p. 41)<sup>1</sup>, people  
have never been as cognitively impoverished as they are today. Major Web platforms such as  
Google and Facebook serve as hubs, distributors, and curators<sup>2</sup>; their algorithms are  
indispensable for navigating the vast digital landscape, and for enabling bottom-up participation  
in the production and distribution of information. Big Tech exploits this all-important role in  
pursuit of the most precious resource in the online marketplace: human attention. Employing  
45 algorithms that learn people's behavioural patterns<sup>4,5,6,7</sup>, technology companies target their users  
with advertisements and design users' information and choice environments<sup>8</sup>. The relationship  
between platforms and people is profoundly asymmetric: Platforms have deep knowledge of  
users' behaviour, whereas users know little about how their data is collected, how it is exploited  
for commercial or political purposes, and how it and the data of others are used to shape their  
50 online experience. These asymmetries in Big Tech's business model have created an opaque  
information ecology that undermines not only user autonomy but also the transparent exchange  
on which democratic societies are built<sup>9,10</sup>. Several problematic social phenomena pervade the  
Internet, such as the spread of false information<sup>11,12,13,14,15</sup>—which includes disinformation  
(intentionally fabricated falsehoods) and misinformation (falsehoods created without intent, e.g.,  
55 poorly researched content or biased reporting)—or attitudinal and emotional polarization<sup>16,17</sup>  
(e.g., polarization of elites<sup>18</sup>, partisan sorting<sup>19</sup>, and polarization along controversial topics<sup>20,21</sup>).  
We argue that the behavioural sciences should play a key role in informing and designing  
systematic responses to such threats. The role of behavioural science is not only to advance  
active scientific debates on the causes and reach of false information<sup>22,23,24,25,26</sup> or whether mass  
60 polarization is increasing<sup>27,28,29</sup>; it is also to find new ways to promote the Internet's potential to  
bolster rather than undermine democratic societies<sup>30</sup>. Many global problems—from climate

Deleted: .....Page Break.....

Deleted: ; they

Deleted: <sup>3</sup>

Deleted: they

Deleted: foster

[change to the coronavirus pandemic—require coordinated collective solutions, making a democratically interconnected world crucial](#)<sup>31</sup>.

### Why Behavioural Sciences Are Crucial for Shaping the Online Ecosystem

More than any traditional media, online media permit and encourage active behaviours<sup>32</sup> such as information search, interaction, and choice. These behaviours are highly contingent on environmental and social structures and cues<sup>33</sup>. Even seemingly minor aspects of the design of digital environments can shape individual actions and scale up to notable changes in collective behaviours. For instance, curtailing the number of times a message can be forwarded on WhatsApp (thereby slowing large cascades of messages) may have been a successful response to the spread of misinformation in Brazil and India<sup>34</sup>.

To a substantial degree, social media and search engines have [taken on a role](#) as [intermediary](#) gatekeepers [between readers and publishers](#): Today, more than half (55%) of global Internet users turn to either social media or search engines to access news [articles](#)<sup>3</sup>. One implication of this seismic shift is that a small number of global corporations and Silicon Valley CEOs have significant responsibility for curating the general population's information<sup>35</sup>—and, [by implication](#), for interpreting and protecting civic freedoms. The flow of information depends on [corporations'](#) willingness and ability to self-regulate the industry. Facebook's recent decision to declare politicians and their ads off-limits to their third-party fact checkers illustrate how corporate decisions can affect citizens' information ecology and the interpretation of fundamental rights, such as freedom of speech. [This situation, in which political content and news diets are curated by opaque and largely unaccountable third parties, is considered unacceptable by a majority of the public](#)<sup>36,37</sup>, who continue to be concerned about their ability to [discern online what is true and what is false](#)<sup>3</sup> and [rate accuracy a very important attribute for social media sharing](#)<sup>38</sup>.

Deleted: <sup>2</sup>

Deleted: replaced traditional media<sup>2</sup>

Deleted: key

Deleted: to news

Formatted: Endnote Reference, Font colour: Auto, German

Deleted: <sup>2</sup>

Deleted: to that extent

Deleted: their

Formatted: Endnote Reference, Font colour: Auto, German

Deleted: <sup>2</sup>

100 How can citizens and democratic governments be empowered<sup>39</sup> to create an ecosystem  
 that “values and promotes truth” (p. 1096)<sup>15</sup>? The answers must be informed by independent  
 behavioural research, which can then form the basis both for improved self-regulation by the  
 relevant companies and for government regulation<sup>40,41</sup>. Regulators in particular face three serious  
 problems in the online domain, that underscore the importance of enlisting the behavioural  
 sciences. The first problem is that online platforms can leverage their proprietary knowledge of  
 105 user behaviour to defang regulations. An example comes from most of the current consent forms  
 under the EU General Data Protection Regulation: Instead of obtaining genuinely informed  
 consent, the current dialogue boxes influence people’s decision-making through self-serving  
 forms of choice architecture (e.g., consent is assumed from pre-ticked boxes or inactivity)<sup>42,43</sup>.  
 This example highlights the need for industry-independent behavioural research in order to  
 110 ensure transparency for the user and to avoid opportunistic responses by those who are regulated.  
 The second problem is that the speed and adaptability of technology and its users exceed that of  
 regulation directly targeting online content. If uninformed by behavioural science, any regulation  
 that focuses only on the symptoms and not on the actual human-platform interaction could be  
 quickly circumvented by users and platforms. The third problem is the risk of censorship  
 115 inherent in regulations that target content; behavioural sciences can reduce that risk as well.  
 Rather than deleting or flagging posts based on judgements about their content, we focus here on  
 how to redesign digital environments so as to provide a better sense of context and to encourage  
 and empower people to make critical decisions for themselves<sup>44,45,46</sup>.

120 Our aim is to enlist two streams of research that illustrate the promise of behavioural  
 sciences. The first examines the informational cues that are available online<sup>32</sup> and asks which can  
 help users gauge the epistemic quality of content or the trustworthiness of the social context from  
 which it originated. The second stream concerns the use of meaningful and predictive cues in

**Deleted:** instead

**Formatted:** Endnote Reference, Font colour: Auto, German

**Deleted:** <sup>14</sup>

**Deleted:** , which

**Deleted:** in order to achieve effective Web governance

**Deleted:** regulation

**Deleted:** forms

**Deleted:** (ideally with enhanced data access)

**Formatted:** Endnote Reference, Font colour: Auto, German

**Deleted:** <sup>30</sup>

behavioural interventions. Interventions can take the form of nudging<sup>47</sup>, which alters the environment or choice architecture so as to draw users' attention to these cues, or boosting<sup>48</sup>, which teaches users to search for them on their own, thereby helping them become more resistant to false information and manipulation in the long run.

### **Digital Cues and Behavioural Interventions for Human-Centred Online Environments**

The online world has the potential to provide digital cues that can help people assess the epistemic quality of content<sup>49,50,51</sup>—the potential of self-contained units of information (here we focus on online articles and social media posts) to contribute to true beliefs, knowledge, and understanding—and the public's attitudes to societal issues<sup>52,53</sup>. We classify those cues as *endogenous* or *exogenous*<sup>54</sup>.

Endogenous cues refer to the content itself, like the plot or the actors and their relations. Modern search engines use natural language-processing tools that analyse content<sup>55</sup>. They have considerable virtues and promise, but current results rarely afford nuanced interpretations<sup>56</sup>. For example, these methods cannot reliably distinguish between facts and opinions, nor can they detect irony, humour, or sarcasm<sup>57</sup>. They also have difficulty differentiating between extremist content and counterextremist messages<sup>58</sup> because both types of messages tend to be tagged with similar keywords. A more general shortcoming of current endogenous cues of epistemic quality is that their evaluation requires background knowledge of the issue in question, which often makes them non-transparent and potentially prone to abuse for censorship purposes.

By contrast, exogenous cues are easier to harness as indicators of epistemic quality: They refer to the context of information rather than the content, are relatively easy to quantify, and can be interpreted intuitively. A famous example of the use of exogenous cues is Google's PageRank algorithm<sup>81</sup>, which takes centrality as a key indicator of quality: Well-connected websites appear

higher up in search results, irrespective of their content. Exogenous cues can indicate how well a piece of information is embedded in existing knowledge or the public discourse.

From here on we focus on exogenous cues and how they can be enlisted by nudging<sup>47</sup> and boosting<sup>48</sup>. Let us emphasize that a single measure will not reach everyone in a heterogeneous population with diverse motives and behaviours. We therefore propose a range of measures that differ in their scope and in the level of user engagement required. Nudging interventions shape behaviour primarily through the design of choice architectures and typically require little active user engagement. Boosting interventions, in contrast, focus on creating and promoting cognitive and motivational competences, either by directly targeting competences as external tools or indirectly by enlisting the choice environment. They require some level of user engagement and motivation. Both nudging and boosting have been shown to be effective in various domains, including health<sup>59,60</sup> and finances<sup>61</sup>. Recent empirical results from research on people's ability to detect false news indicate that informational literacy can also be boosted<sup>62</sup>. Initial results on the effectiveness of simple nudging interventions that remind people to think about accuracy before sharing content<sup>38</sup> also suggest that interventions based on behavioural sciences could be effective in the online domain<sup>64</sup>. While empirical tests and evidence are urgently needed, the first step is to outline the conceptual space of possible interventions and make specific proposals.

Table 1 examines three online contexts: articles from newspapers or blogs, algorithmic curation systems that automatically suggest products or information (e.g., search engines or algorithmic curation of news feeds), and social media that display information about the behaviour of others (e.g., shared posts or social reactions such as comments or "likes"). Each is associated with a unique set of challenges, cues, and potential interventions. Next, we review the challenges and cues in Table 1, and detail some interventions in the subsequent sections.

**Deleted:** two classes of behavioural interventions:

**Deleted:** <sup>42</sup>

**Formatted:** Endnote Reference, Font colour: Auto, German

**Formatted:** Endnote Reference, Font colour: Auto, German

**Deleted:** <sup>43</sup>

**Deleted:** A

**Deleted:** <sup>63</sup>

**Formatted:** Indent: First line: 0 cm

Context	Challenges	Cues	Nudging	Boosting
Online articles	Information overload and fragmentation of sources	Cues to epistemic quality, like references	...to pay attention to epistemic cues and external evidence.	...routines to systematically check epistemic cues.
Algorithmic curation	Asymmetry of knowledge and opaque manipulation	Transparent recommendation and sorting criteria	...awareness of factors that shape recommendations and the news feed.	...self-nudging towards quality information.
Social media	Lack of global network information false consensus effects	Global social cues that include base rates and passive behaviour	...to consider global social cues and accuracy before sharing.	...to infer credibility from social context and history of content.

Deleted:

Deleted: <object>

**Table 1.** Overview of challenges, cues, and potential targets of nudging and boosting interventions in three online contexts.

**Online Articles: Information Overload and Epistemic Cues.** The capacity to transfer information online continues to increase exponentially (average annual growth rate: 28%)<sup>65</sup>. Content can be distributed more rapidly and reaches an audience faster<sup>66</sup>. This increasing pace has consequences. In 2013, a hashtag on Twitter remained in the top 50 most popular hashtags worldwide for an average of 17.5 hours; by 2016, a hashtag's time in the limelight had dropped to 11.9 hours. The same declining half-life has been observed for Google queries and movie ticket sales<sup>67</sup>. This acceleration, arguably driven by the [finite](#) limits of attention available for the ever increasing quantity of topics and content<sup>68</sup> combined with an apparent thirst for novelty has significant but underappreciated psychological consequences. Information overload makes it harder for people to make good decisions about what to look at, spend time on, believe, and share<sup>69,70</sup>. For instance, longer-term offline decisions such as choosing a newspaper subscription (that then constrains one's information diet) have evolved into a multitude of online micro-decisions about which individual articles to read from a scattered array of sources. The more

Deleted: invariant



210 sources crowd the market, the less attention can be allocated to each piece of content, and the  
more difficult it becomes to assess their trustworthiness— even more so given the demise [and](#)  
[erosion](#) of classic indicators of quality<sup>71</sup> (e.g., name recognition, reputation, print quality, price).  
[Going beyond them, new](#) cues for epistemic quality that are readily accessible even under  
information overload are necessary. Exogenous cues can highlight the epistemic quality of  
215 individual articles, in particular by showing how an article is embedded in the existing corpus of  
knowledge and public discourse. These cues include, for instance, a newspaper article’s sources  
and citation network (i.e., sources that cite the article or are cited by it), references to established  
concepts and topical empirical evidence, and even the objectivity of the language.

**Algorithmic Curation: Asymmetry of Knowledge and Transparency.** To help users  
220 navigate the overabundance of information, search engines automatically order results<sup>72,73</sup> and  
recommender systems<sup>74</sup> guide users to content they are likely to prefer<sup>75</sup>. But this convenience  
[exact](#)s a price. Because user satisfaction is not necessarily in line with the goals of algorithms—  
to maximize user engagement and screen time<sup>76</sup>—algorithmic curation often deprives users of  
autonomy. For instance, feedback loops are created that can artificially [re-enforce](#) preferences  
225 <sup>77,78,79,80</sup>, and recommender systems can eliminate context in order to avoid overburdening users.

[To stay up to date and engaging](#), algorithms can trade recency for importance<sup>81</sup> [and, by](#)  
[optimizing on click rates, trade](#) “clickbait” for quality.

Deleted: New

Deleted: comes at

Deleted: amplify

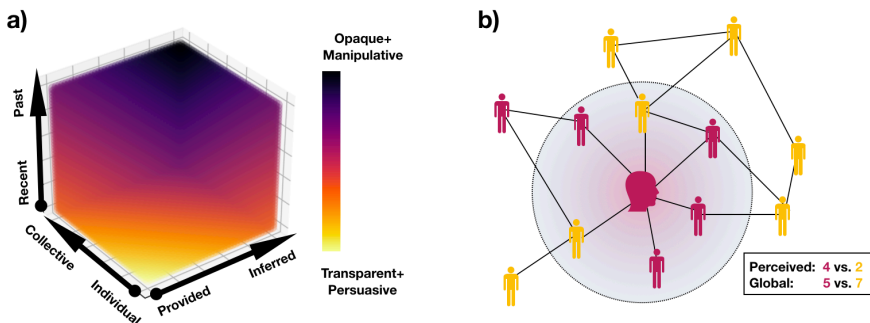
Deleted: By optimizing on click rates,

Deleted: and

Deleted: <sup>82</sup>

Similarly, aggregated previous user selections make targeted commercial nudging—and even manipulation—possible<sup>83,84</sup>. For example, given just 300 Facebook “likes” from one person, a regression model can better predict that person’s personality traits than friends and family<sup>85</sup>. There are at least three dimensions of knowledge where platforms can far exceed individual human capabilities (Figure 1a): Data that reaches further back in time (e.g., years of location history on Google Maps), information about behaviour on a collective rather than an individual level (e.g., millions of Amazon customers with similar interests [as an individual can be utilized to recommend further products to that person](#)), and knowledge that is inferred from existing data using machine-learning methods ([e.g. food preferences from movement patterns between restaurants](#)).

Moving further along these dimensions, it becomes more difficult for a user to comprehend the wealth and predictive potential of this knowledge. Automatic customization of online environments that is based on this knowledge can therefore be opaque and manipulative (Figure 1a). [Recent surveys in the USA and Germany found that a majority of respondents consider such data-driven personalization of political content \(61%\), social media feeds \(57%\) and news diets \(51%\) unacceptable, whereas they are much more accepting of it when it pertains](#)



**Figure 1.** Challenges in automatically curated environments and on social media platforms. (a) Dimensions of knowledge that platforms can acquire with information technology, which make their recommendations continuously opaque and manipulative. (b) perceived group sizes versus the actual global sizes, from the viewpoint of one user (head icon in the center) in a homophilic social network.

[to commercial content](#). To rebalance the relationship between algorithmic and human decision making [and to allow for heterogeneous preferences across different domains](#), a two-step process is required. First, steps should be taken toward [the design and implementation of](#) more transparent algorithms. They should offer cues that clearly represent the data types and the weighting that led to a system's suggestion as well as offer information about the target audience. Second, users should be able to adapt these factors to their personal preferences in order to regain autonomy.

**Social Media: Network Effects and Social Cues.** More than two thirds of all Internet users, around 3 billion people, actively use social media<sup>86</sup>. These platforms offer information about the behaviour of others (e.g., “likes” and emoticons)<sup>87</sup> and new opportunities for interaction (e.g., follower relationships and comment sections). However, these signals and interactions are often one-dimensional, represent only a user's immediate online neighbourhood, and do not distinguish between different types of connections<sup>88</sup>. These limitations can have drastic effects, such as dramatically changing a user's perception of group sizes<sup>89,90</sup> and giving rise to false-consensus effects (i.e., the majority opinion in an individual's neighbourhood leads them to falsely believe it reflects the actual majority opinion; Figure 1b). When people associate with like-minded others from a globally dispersed online community, their self-selected social surroundings (known as a homophilic social network) and the low visibility of the global state of the network<sup>91,92</sup> can create the illusion of broad support<sup>93</sup> and reinforce opinions or even make them more extreme<sup>94,95</sup>. For instance, even if only a tiny fraction (e.g., one in a million) of the more than two billion Facebook users believe that the Earth is flat, they could still form an online community of thousands, thereby creating a shield of like-minded people against corrective efforts<sup>96,97,98,99</sup>. Although large social media platforms routinely aggregate information that would foster a realistic assessment of societal attitudes, they currently do not provide a well-calibrated impression of the degree of public consensus<sup>100</sup>. Instead, they show reactions from

275 others as asymmetrically positive—there typically is no “dislike” button—or biased toward  
narrow groups or highly active users<sup>101</sup> in order to maximize user engagement. This need not be  
the case: The interactive nature of social media could be harnessed to promote diverse  
democratic dialogue and foster collective intelligence. In order to achieve this goal, social media  
needs to offer more meaningful, higher-dimensional cues that carry information about the  
280 broader state of the network rather than just the user’s direct neighbourhood, which can mitigate  
biased perceptions caused by the network structure<sup>102</sup>. For instance, social media platforms could  
provide a transparent crowd-sourced voting system<sup>103</sup> or display informative metrics about the  
behaviour and reactions of others (e.g., including passive behaviour, like the total number of  
people who scrolled over a post), which might counter false-consensus effects.

285 **Nudging Interventions to Shape Online Environments**

Nudging interventions can alter choice architectures to promote the epistemic quality of  
information and its spread. One type of nudge, educative nudging, integrates epistemic cues into  
the choice environment primarily to inform behaviour (as opposed to actively steering it). For  
instance, highlighting when content stems from few or anonymous sources (as used by

Deleted: .

Deleted: <object>



**Figure 2** Nudging interventions that modify online environments. (a) Examples of exogenous cues and how they could appear alongside a social media post. (b) Example of a transparently organized news feed on social media. Types of content are clearly distinguished, sorting criteria and their values are shown with every post, and users can adjust weightings.

Wikipedia) can remind people to scrutinize content more thoroughly<sup>104,105</sup> and simultaneously create an incentive structure for content producers to meet the required criteria. Such outlets can be made more transparent, for example by disclosing the identity of their confirmed owners.

295 Similarly, pages that are run by state-controlled media might be labelled as such<sup>106</sup>. Going a step further, adding prominent hyperlinks to vetted reference sources for important concepts in a text could encourage a reader to gain context by perusing multiple sources—a strategy used by professional fact checkers<sup>107</sup>.

300 Nudges can also communicate additional information about what others are doing, thereby invoking the steering power of descriptive social norms<sup>108</sup>: For instance, contextualizing the number of likes [by expressing them against the absolute frequency of total readers](#) (e.g., “4,287 [out of 1.5 million](#) readers liked this article”) might counteract false-consensus effects that a number presented without context (“4,287 people liked this article”) may otherwise engender. Transparent numerical formats have already been shown to successfully improve statistical

305 literacy in the medical domain<sup>109</sup>. Similarly, displaying the total number of readers and their average reading time in relation to the potential total readership could help users evaluate [the](#) content’s epistemic quality: If only a tiny portion of the potential readership has actually read an article, [whereas the majority spent less than a second on it](#), it might be clickbait. The presentation of many other cues, including ones that reach into the history of a piece of content, could be used

310 to promote epistemic value on social media. Figure 2a shows a nudging intervention that integrates several exogenous cues into a social media news feed.

Similarly, users can be discouraged from sharing low-quality information without resorting to censorship by introducing “friction”—for instance, by making the act of sharing slightly more effortful<sup>110</sup>. In this case, sharing low-quality content may require a further mouse

Deleted: ,502,961

click [in a pop-up warning message](#), alongside additional information [about which of the above](#)  
cues [are missing or have critical values](#).

Deleted: based on

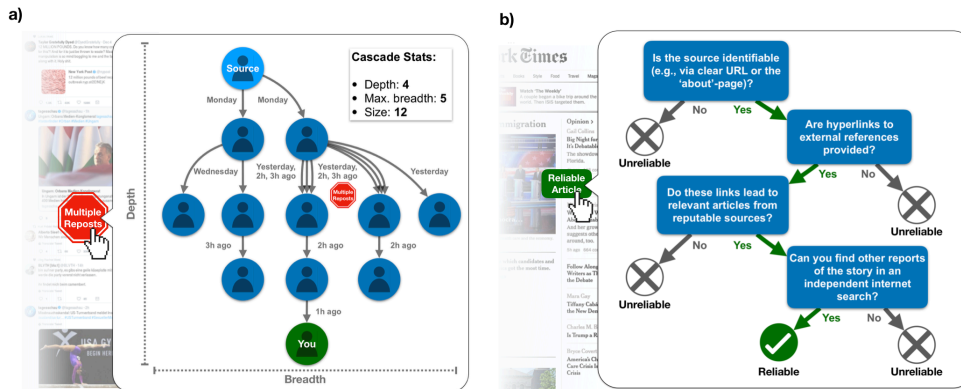
Deleted: discussed above.

Deleted: <object>

Another type of nudge targets how content is arranged in browsers. The way a social  
media news feed sorts content is crucial in shaping how much attention is devoted to particular  
posts. Indeed, news feeds have become one of the most sophisticated algorithmically driven  
choice architectures of online platforms<sup>8,111</sup>. Transparent sorting algorithms for news feeds (such  
as the algorithm used by Reddit) that show the factors that determine how posts are sorted can  
help people understand why they see certain content; at the very least this nudging intervention  
would make the design of the feed's architecture more transparent. Relatedly, platforms that  
clearly differentiate between types of content (e.g., ads, news, or posts by friends) can make  
news feeds more transparent and clearer (Figure 2b).

Formatted: Endnote Reference, Font colour: Auto, German

Deleted: <sup>7</sup>



**Figure 3** Illustrations of boosting interventions as they could appear within an online environment or as external tools. (a) Visualization of a sharing cascade. Alongside metrics, like the depth or the breadth of a cascades, a pop-up window on social media can provide a simple visualization of a sharing cascade that shows who (if the profile is public) and when others have shared content before it reached the user. (b) A fast-and-frugal decision tree as an example of a boosting intervention. A pop-up or an external tool can show a fast-and-frugal decision tree alongside an online article, that helps to check criteria to evaluate its reliability, where the criteria were adapted from professional fact checkers and primarily point to checking external information<sup>90</sup>.

### Boosting interventions to foster user competences

Boosting seeks to empower people in the longer term by helping them build the competences they need to navigate situations autonomously (for a conceptual map of boosting interventions online, see also<sup>112</sup>). These interventions can be integrated directly into the environment itself or be available in an app or browser add-on. Unlike many nudging interventions, boosting interventions will ideally remain effective even when they are no longer present in the environment because they have become routinized and have instilled a lasting competence in the user.

The competence of acting as one's own choice architect, or self-nudging, can be boosted<sup>113</sup>. For instance, when users can customize how their news feed is designed and sorted (Figure 2b), they can become their own choice architects and regain some informational autonomy. For instance, users could be enabled or encouraged to design information ecologies for themselves that are tailored toward high epistemic quality, making sources of low epistemic quality less accessible. Such boosting interventions would require changes to the online environment (e.g., transparent sorting algorithms or clear layouts; see previous section and Figure 2b) and the provision of epistemic cues.

Another competence that could be boosted to help users deal more expertly with information they encounter online is the ability to make inferences about the reliability of information based on the social context from which it originates<sup>114</sup>. The structure and details of the entire cascade of individuals who have previously shared an article on social media has been shown to serve as proxies for epistemic quality<sup>115</sup>. [Namely, the](#) sharing cascade contains metrics such as the depth and breadth of dissemination by others, with deep and narrow cascades indicating extreme or niche topics and breadth indicating widely discussed issues<sup>116</sup>. A boosting intervention could provide this information (Figure 3a) to show the full history of a post, including the original source, the friends and public users who disseminated it, and the timing of

Deleted: The

the process (showing, e.g., if the information is old news that has been repeatedly and artificially amplified).

Yet another competence required for discerning between sources of high and low quality is the ability to read laterally.<sup>107</sup> Lateral reading is a skill developed by professional fact checkers that entails looking for information on sites other than the information source in order to evaluate its credibility (e.g., “who is behind this website?” and “what is the evidence for its claims?”)

rather than evaluating a website’s credibility by using the information provided there. This competence can be boosted with simple decision aids such as fast-and-frugal decision trees<sup>117,118</sup>.

Employed in a wide range of areas (e.g., medicine, finance, law, management), fast-and-frugal decision trees can guide the user to scrutinize relevant cues. For example, users can respond to prompts in a pop-up window (e.g., “Are references provided?”), with each answer leading either to an immediate decision (e.g., “unreliable”) or to the next cue until a final judgment about content reliability is reached (e.g., “reliable”; Figure 3b)<sup>119</sup>. Decision trees can also enhance the

transparency of third-party decisions. If reliability is judged by third-party fact checkers or via an automated process, users could opt to see the decision tree and follow the path that led to the decision, thereby gaining insight that will be useful in the long-term. Eventually, fast-and-frugal decision trees may help people establish a habit of checking epistemic cues when reading content even in the absence of a pop-up window suggesting they do so.<sup>48</sup>

Finally, the competence of understanding what makes intentionally false information so alluring (e.g., novelty and the element of surprise) can be boosted by mental inoculation techniques: Being informed about manipulative methods before encountering them online enables an individual to detect parasitic imitations of trustworthy sources and other sinister tactics<sup>120,121</sup>. Making people aware of such strategies or of their own personal vulnerabilities leaves them better able to identify and resist manipulation. For instance, having people take on

Deleted: <sup>100</sup>

Formatted: Endnote Reference, Font colour: Auto, German

Deleted: it provides

Deleted: to

Formatted: Endnote Reference, Font colour: Auto, German

Deleted: <sup>43</sup>



385 the role of a malicious influencer in a computer game has been demonstrated to improve their  
ability to spot and resist misinformation<sup>62,122</sup>. This inoculation technique can be used in a range  
of contexts online; for example, learning about the target group of an advertisement can increase  
people's ability to detect advertising strategies.

### Conclusion

390 Any attempt to regulate or manage the digital world must begin with the understanding  
that online communication is already regulated— to some extent by public policy or laws, but  
primarily by search engines and recommender systems whose goals and parameters may not be  
publicly known, let alone be subject to public scrutiny. The current online environment has given  
rise to opaque and asymmetric relationships between users and platforms, and it is reasonable to  
395 question whether the industry will independently take action to foster an ecosystem that values  
and promotes truth. The interventions we propose are aimed primarily at empowering individuals  
to make informed and autonomous decisions in the online ecosystem—and, through their own  
behaviour, to foster and reinforce truth. The interventions are partly conceptualized on the basis  
of existing results. However, not all interventions have been tested in the specific context in  
400 which they may be deployed. Undoubtedly, therefore, these and other interventions are subject to  
further empirical test. The first results are promising, identifying some interventions as  
effective<sup>63,121</sup> whereas others appear less promising<sup>123</sup>. Both set of results will inform the design  
of more effective interventions.

In our view, the future task for scientists is to design interventions that meet at least three  
405 selection criteria: They must be transparent and trustworthy to the public, standardisable within  
certain categories of content, and, importantly, hard to game by bad-faith actors or vested  
interests. We also suggest that there is a need to examine a wide spectrum of interventions, from  
nudges to boosts, in order to reach different types of people, who have heterogeneous

Deleted: <sup>57</sup>

Formatted: Endnote Reference, Font colour: Auto, German

Deleted: also

Deleted: The

Deleted: <sup>62</sup>

Deleted: <sup>119</sup>

Deleted: . These interventions will not entirely

415 [preferences, motivations and online behaviours. These interventions will not completely](#) prevent  
manipulation or active dissemination of false information, but they will help users recognise  
when these malicious tactics are at work. They will also permit producers of quality information  
to differentiate themselves from less trustworthy sources. Behavioural interventions in the online  
ecology can not only inform government regulations, but also signal a platform's commitment to  
420 truth, epistemic quality, and trustworthiness: Platforms can indicate their commitment to these  
values by providing their users with exogenous cues and boosting and nudging interventions, and  
users can choose to avoid platforms that do not offer them these features.

For this dynamic to gain momentum it is not necessary that all or even the majority of  
users engage with nudging or boosting interventions; as the first Wikipedia contributors have  
425 proven, a critical mass may suffice to allow positive effects to scale up to major improvements.  
[Such a dynamic may counteract a possible drawback of the proposed interventions; namely,  
widening information gaps between users if only empowered consumers are able to recognise  
quality information. Furthermore, it can help to mitigate gaps currently arising from the ability to  
pay for quality content. In the trade-off between offering interventions that not everybody will  
430 engage with and leaving citizen without any interventions that could cause differential  
competences and knowledge, we err on the side of empowerment.](#)

**Acknowledgments:** We thank Anastasia Kozyreva and Stefan Herzog for their helpful  
comments, and Deb Ain for editing the manuscript. Ralph Hertwig and Stephan Lewandowsky  
435 acknowledge support from the Volkswagen Foundation. The funders had no role in study design,  
data collection and analysis, decision to publish or preparation of the manuscript. **Competing  
interests:** Cass R. Sunstein has [served as a paid consultant on two occasions](#) for Facebook.

Formatted: Indent: First line: 1.25 cm, Pattern: Clear

Deleted: occasionally

Deleted: an advisor

**Author contributions:** P.L.S., S.L. and R.H. conceptualized the project, all authors wrote the manuscript.

## References:

<sup>1</sup> Simon H.A., Designing organizations for an information-rich world. *Computers, Communications and the Public Interest* **70**, 37-72 (John Hopkins Press, 1971).

<sup>2</sup> Newman, N., Fletcher, R., Kalogeropoulos, A., & Nielsen, R. *Reuters institute digital news report 2019* (Reuters Institute for the Study of Journalism, 2019).

<sup>4</sup> Kosinski M., Stillwell D., & Graepel T., Private traits and attributes are predictable from digital records of human behavior. *Proc. Natl. Acad. Sci. U.S.A* **110**, 5802-5805 (2012).

<sup>5</sup> Boerman, S. C., Kruikemeier, S., & Zuiderveen Borgesius, F. J. Online behavioral advertising: A literature review and research agenda. *Journal of Advertising*, **46**, 363-376 (2017).

<sup>6</sup> Ruths, D., & Pfeffer, J. Social media for large studies of behavior. *Science*, **346**, 1063-1064 (2014).

<sup>7</sup> Tufekci, Z. Engineering the public: Big data, surveillance and computational politics. *First Monday*, **19**, (2014).

<sup>8</sup> Harris, T. How Technology is Hijacking Your Mind—from a Magician and Google Design Ethicist. *Thrive Global*, **18**, (2016).

<sup>9</sup> Persily, N. The 2016 US Election: Can democracy survive the internet? *Journal of Democracy*, **28**, 63-76 (2017).

<sup>10</sup> Habermas, J., *The Structural Transformation of the Public Sphere: An Inquiry Into a Category of Bourgeois Society*. (MIT Press, 1991).

<sup>11</sup> Vosoughi S., Roy D., & Aral S., The spread of true and false news online. *Science* **359**, 1146-1151 (2018).

<sup>12</sup> Mocanu, D., Rossi, L., Zhang, Q., Karsai, M., & Quattrociocchi, W. Collective attention in the age of (mis) information. *Computers in Human Behavior*, **51**, 1198-1204 (2015).

<sup>13</sup> Rich, M. D. *Truth decay: An initial exploration of the diminishing role of facts and analysis in American public life*. (Rand Corporation, 2018).

**Deleted:** <sup>3</sup> Newman, N., Fletcher, R., Kalogeropoulos, A., & Nielsen, R. *Reuters institute digital news report 2019* (Reuters Institute for the Study of Journalism, 2019).

<sup>14</sup> Vargo, C. J., Guo, L., & Amazeen, M. A. The agenda-setting power of fake news: A big data analysis of the online media landscape from 2014 to 2016. *New Media & Society*, **20**, 2028-2049 (2018).

<sup>15</sup> Lazer D. M. J., et al., The science of fake news. *Science* **359**, 1094-1096 (2018).

<sup>16</sup> Baldassarri, D., & Gelman, A., Partisans without constraint: Political polarization and trends in American public opinion. *American Journal of Sociology*, **114**, 408-446 (2008).

<sup>17</sup> Abramowitz, A. I., & Saunders, K. L., Is polarization a myth?. *The Journal of Politics*, **70**, 542-555 (2008).

<sup>18</sup> McCarty, N., Poole, K. T., & Rosenthal, H., *Polarized America: The dance of ideology and unequal riches*. (MIT Press, 2006).

<sup>19</sup> Fiorina, M. P., & Abrams, S. J., Political polarization in the American public. *Annu. Rev. Polit. Sci.*, **11**, 563-588 (2008).

<sup>20</sup> McCright, A. M., & Dunlap, R. E., The politicization of climate change and polarization in the American public's views of global warming, 2001–2010. *The Sociological Quarterly*, **52**, 155-194 (2011).

<sup>21</sup> Cota, W., Ferreira, S.C., Pastor-Satorras, R. *et al.* Quantifying echo chamber effects in information spreading over political communication networks. *EPJ Data Sci.* **8**, 35 (2019)

<sup>22</sup> DiMaggio, P., Evans, J., & Bryson, B., Have American's social attitudes become more polarized?. *American journal of Sociology*, **102**, 690-755 (1996).

<sup>23</sup> Fletcher, R., Cornia, A., Graves, L., & Nielsen, R. K., Measuring the reach of “fake news” and online disinformation in Europe. *Reuters institute factsheet*. (2018).

<sup>24</sup> Cinelli, M., Cresci, S., Galeazzi, A., Quattrociocchi, W., & Tesconi, M. The Limited Reach of Fake News on Twitter during 2019 European Elections. Preprint at: <https://arxiv.org/abs/1911.12039> (2020).

<sup>25</sup> Guess, A. M., Nyhan, B., & Reifler, J., [Exposure to untrustworthy websites in the 2016 US election](#). *Nature human behaviour*, **1-9**, (2020).

<sup>26</sup> Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A., & Bonneau, R., Tweeting from left to right: Is online political communication more than an echo chamber?. *Psychological science*, **26**, 1531-1542 (2015).

<sup>27</sup> Evans, J. H., Have Americans' attitudes become more polarized?—An update. *Social Science Quarterly*, **84**, 71-90 (2003).

**Deleted:** Selective exposure

**Deleted:** misinformation: Evidence from the consumption of fake news during

**Deleted:** presidential campaign. *European Research Council*, ...

**Formatted:** Font: Not Bold

**Deleted:** , (2018)

<sup>28</sup> Lelkes, Y. Mass polarization: Manifestations and measurements. *Public Opinion Quarterly*, **80**, 392-410 (2016).

<sup>29</sup> Del Vicario, M., et al., The spreading of misinformation online. *Proc. Natl. Acad. Sci. U.S.A.*, **113**, 554-559 (2016).

<sup>30</sup> [Watts D. J., Should Social Science Be More Solution-Oriented? \*Nature Human Behaviour\*. \*\*1\*\*, 0015 \(2017\).](#)

<sup>31</sup> [Larson, H. J., The biggest pandemic risk? Viral misinformation, \*Nature\*, \*\*562\*\*, 309-310 \(2018\).](#)

<sup>32</sup> Sundar, S. The MAIN Model : A Heuristic Approach to Understanding Technology Effects on Credibility. *MacArthur Foundation Digital Media and Learning Initiative*. (2007).

<sup>33</sup> Gigerenzer, G., Hertwig, R., & Pachur, T., *Heuristics: The Foundations of Adaptive Behavior*. (Oxford University Press, 2011).

<sup>34</sup> de Freitas Melo, P., Vieira, C. C., Garimella, K., de Melo, P. O. V., & Benevenuto, F., Can WhatsApp Counter Misinformation by Limiting Message Forwarding?. In *International Conference on Complex Networks and Their Applications* (pp. 372-384) (Springer, Cham, 2019).

<sup>35</sup> <https://www.adl.org/news/article/sacha-baron-cohens-keynote-address-at-adls-2019-never-is-now-summit-on-anti-semitism> (accessed: 07.12.2019)

<sup>36</sup> [Kozyreva, A., Herzog, S., Lorenz-Spreen, P., Hertwig, R., & Lewandowsky, S., Artificial intelligence in online environments: Representative survey of public attitudes in Germany. Berlin: Max Planck Institute for Human Development. \(2020\).](#)

<sup>37</sup> [Smith, A., Public attitudes toward computer algorithms. \*Pew Research Center\*. \(2018\).](#)

<sup>38</sup> [Pennycook, G., Epstein, Z., Mosleh, M., Arechar, A. A., Eckles, D., & Rand, D. G., Understanding and reducing the spread of misinformation online. Preprint at: <https://psyarxiv.com/3n9u8/> \(2019\).](#)

<sup>39</sup> S. Zuboff, Surveillance Capitalism and the Challenge of Collective Action. *New Labor Forum*, **28**, 10-29 (2019).

<sup>40</sup> Klein, D., & Wueller, J. Fake news: A legal perspective. *Journal of Internet Law* (2017).

<sup>41</sup> Assemblée nationale, Proposition de loi relative à la lutte contre la manipulation de l'information. (2018). available at: <http://www.assemblee-nationale.fr/15/ta/tap0190.pdf> (accessed: 26.06.2019).

<sup>42</sup> van Ooijen, I., & Vrabec, H. U., Does the GDPR enhance consumers' control over personal data? An analysis from a behavioural perspective. *Journal of Consumer Policy*, **42**, 91-107 (2019).

Deleted: Watts D

Formatted: None

Deleted: Should Social Science Be More Solution-Oriented?

Deleted: *Human Behaviour*. **1**, 0015 (2017)

Formatted: English (US)

Formatted: English (US)

Formatted: English (US)

Formatted: Font: Cambria

- 
- <sup>43</sup> Nouwens, M., Liccardi, I., Veale, M., Karger, D., & Kagal, L., Dark Patterns after the GDPR: Scraping Consent Pop-ups and Demonstrating their Influence. Preprint at: <https://arxiv.org/abs/2001.02479> (2020).
- <sup>44</sup> Hertwig, R., When to consider boosting: some rules for policy-makers. *Behavioural Public Policy*, **1**, 143-161 (2017).
- <sup>45</sup> Epstein, Z., Pennycook, G., & Rand, D. Letting the crowd steer the algorithm: Laypeople can effectively identify misinformation sources. Preprint at: <https://psyarxiv.com/z3s5k/> (2019).
- <sup>46</sup> Britt, M. A., Rouet, J. F., Blaum, D., & Millis, K., A Reasoned Approach to Dealing With Fake News. *Policy Insights from the Behavioral and Brain Sciences*, **6**, 94-101 (2019).
- <sup>47</sup> Thaler R. H., Sunstein. C. R., *Nudge: Improving Decisions About Health, Wealth, and Happiness* (Yale University Press, New Haven, CT, 2008)
- <sup>48</sup> Hertwig R., Grüne-Yanoff T., Nudging and boosting: steering or empowering good decisions. *Perspect Psychol Sci* **12**, 973–986 (2017).
- <sup>49</sup> Griffiths, K. M., & Christensen, H. Website quality indicators for consumers. *Journal of medical Internet research*, **7**, e55 (2005).
- <sup>50</sup> Nickel, M., Murphy, K., Tresp, V., & Gabrilovich, E., A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, **104**, 11-33 (2015).
- <sup>51</sup> Dong, X., et al., Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 601-610 ACM (2014).
- <sup>52</sup> Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. Fake news detection on social media: A data mining perspective. *ACM SIGKDD Explorations Newsletter*, **19**, 22-36 (2017).
- <sup>53</sup> Klačnja, M., Barberá, P., Beauchamp, N., Nagler, J., & Tucker, J., Measuring public opinion with social media data. In *The Oxford handbook of polling and survey methods*. (2017).
- <sup>54</sup> Dong, X. L., et al., Knowledge-based trust: Estimating the trustworthiness of web sources. *Proceedings of the VLDB Endowment*, **8**, 938-949 (2015).
- <sup>55</sup> Jeremy, H., *Google Hummingbird: where no search has gone before*. *Wired*. available at: <https://www.wired.com/insights/2013/10/google-hummingbird-where-no-search-has-gone-before/> (accessed: 09.07.2019)
- <sup>56</sup> Luo, H., Liu, Z., Luan, H., & Sun, M., Online learning of interpretable word embeddings. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 1687-1692 (2015).

<sup>57</sup> Schmidt, A., & Wiegand, M., A survey on hate speech detection using natural language processing. In *Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media*, 1-10 (2017).

<sup>58</sup> Schmitt, J. B., Rieger, D., Rutkowski, O., & Ernst, J. Counter-messages as Prevention or Promotion of Extremism?! The Potential Role of YouTube: Recommendation Algorithms. *Journal of Communication*, **68**, 780-808 (2018).

<sup>59</sup> Arno, A., & Thomas, S. The efficacy of nudge theory strategies in influencing adult dietary behaviour: a systematic review and meta-analysis. *BMC public health*, **16**, 676 (2016).

<sup>60</sup> Kurvers, R. H., et al., Boosting medical diagnostics by pooling independent judgments. *Proc. Natl. Acad. Sci. U.S.A.*, **113**, 8777-8782 (2016).

<sup>61</sup> Lusardi, A., & Mitchell, O. S., The economic importance of financial literacy: Theory and evidence. *Journal of economic literature*, **52**, 5-44 (2014).

<sup>62</sup> Roozenbeek, J., & van der Linden, S., Fake news game confers psychological resistance against online misinformation. *Palgrave Communications*, **5**, (2019).

<sup>64</sup> Pennycook, G., & Rand, D. G., Lazy, not biased: Susceptibility to partisan fake news is better explained by lack of reasoning than by motivated reasoning. *Cognition*, **188**, 39-50 (2019).

<sup>65</sup> Hilbert M., & López P., The world's technological capacity to store, communicate, and compute information. *Science* **332**, 60-65 (2011).

<sup>66</sup> Rosa, H., *Social Acceleration: A New Theory of Modernity*. (New York: Columbia University Press, 2013).

<sup>67</sup> Lorenz-Spreen, P., Mønsted, B. M., Hövel, P., & Lehmann, S., Accelerating Dynamics of Collective Attention. *Nature Communications* **10**, 1759 (2019).

<sup>68</sup> Wu, F., & Huberman, B. A., Novelty and collective attention. *Proc. Natl. Acad. Sci. U.S.A.*, **104**, 17599-17601 (2007).

<sup>69</sup> Hills, T. T., Noguchi, T., & Gibbert, M., Information overload or search-amplified risk? Set size and order effects on decisions from experience. *Psychonomic Bulletin & Review*, **20**, 1023-1031 (2013).

<sup>70</sup> Hills, T. T., The Dark Side of Information Proliferation. *Perspectives on Psychological Science*, **14**, 323–330 (2019).

<sup>71</sup> American Society of News Editors (ASNE), Statement of principles. Available at: <https://www.asne.org/content.asp?pl=24&sl=171&contentid=171> (accessed: 27.05.2019).

**Deleted:** <sup>63</sup> Pennycook, G., Epstein, Z., Mosleh, M., Arechar, A. A., Eckles, D., & Rand, D. G., Understanding and reducing the spread of misinformation online. Preprint at: <https://psyarxiv.com/3n9u8/> (2019).<sup>†</sup>

<sup>72</sup> Epstein, R., & Robertson, R. E., The search engine manipulation effect (SEME) and its possible impact on the outcomes of elections. *Proc. Natl. Acad. Sci. U.S.A.*, **112**, E4512-E4521 (2015).

<sup>73</sup> Lazer, D., The rise of the social algorithm. *Science*, **348**, 1090-1091 (2015).

<sup>74</sup> Resnick, P., & Varian, H. R., Recommender systems. *Communications of the ACM*, **40**, 56-59 (1997).

<sup>75</sup> Bakshy, E., Messing, S., & Adamic, L. A., Exposure to ideologically diverse news and opinion on Facebook. *Science*, **348**, 1130-1132 (2015).

<sup>76</sup> Martens, Be., Aguiar, L., Gomez-Herrera, E., Mueller-Langer, F., The Digital Transformation of News Media and the Rise of Disinformation and Fake News. Digital Economy Working Paper 2018-02, *Joint Research Centre Technical Reports*. Available at: <https://ssrn.com/abstract=3164170> (2018).

<sup>77</sup> Cosley, D., Lam, S. K., Albert, I., Konstan, J. A., & Riedl, J., Is seeing believing?: how recommender system interfaces affect users' opinions. In *Proceedings of the SIGCHI conference on Human factors in computing systems* 585-592 ACM (2003).

<sup>78</sup> Pan, B., Hembrooke, H., Joachims, T., Lorigo, L., Gay, G., & Granka, L., In google we trust: Users' decisions on rank, position, and relevance. *Journal of Computer-Mediated Communication*, **12**, 801-823 (2007).

<sup>79</sup> Bozdag, E., Bias in algorithmic filtering and personalization. *Ethics and Information Technology*, **15**, 209-227 (2013).

<sup>80</sup> Sunstein, C. R., *Republic. com*. (Princeton university press, 2002).

<sup>81</sup> Chakraborty, A., Ghosh, S., Ganguly, N., & Gummadi, K. P., Optimizing the recency-relevancy trade-off in online news recommendations. In *Proceedings of the 26th International Conference on World Wide Web* 837-846 (2017).

<sup>83</sup> Zuboff, S., Big other: surveillance capitalism and the prospects of an information civilization. *Journal of Information Technology*, **30**, 75-89 (2015).

<sup>84</sup> Matz, S. C., Kosinski, M., Nave, G., & Stillwell, D. J., Psychological targeting as an effective approach to digital mass persuasion. *Proc. Natl. Acad. Sci. U.S.A.*, **114**, 12714-12719 (2017).

<sup>85</sup> Youyou, W., Kosinski, M., & Stillwell, D., Computer-based personality judgments are more accurate than those made by humans. *Proc. Natl. Acad. Sci. U.S.A.*, **112**, 1036-1040 (2015).

<sup>86</sup> <https://ourworldindata.org/rise-of-social-media> (accessed: 05.12.2019)

**Deleted:** <sup>82</sup> Chakraborty, A., Ghosh, S., Ganguly, N., & Gummadi, K. P., Optimizing the recency-relevancy trade-off in online news recommendations. In *Proceedings of the 26th International Conference on World Wide Web* 837-846 (2017).¶



- 
- <sup>87</sup> Porten-Che  , P., & Eilders, C., The effects of likes on public opinion perception and personal opinion. *Communications, The European Journal of Communication Research* (2019)
- <sup>88</sup> Dandekar, P., Goel, A., & Lee, D. T., Biased assimilation, homophily, and the dynamics of polarization. *Proc. Natl. Acad. Sci. U.S.A.*, **110**, 5791-5796 (2013).
- <sup>89</sup> Lee, E., Karimi, F., Wagner, C., Jo, H. H., Strohmaier, M., & Galesic, M., Homophily and minority-group size explain perception biases in social networks. *Nature human behaviour*, **3**, 1078-1087 (2019).
- <sup>90</sup> Stewart, A. J., Mosleh, M., Diakonova, M., Arechar, A. A., Rand, D. G., & Plotkin, J. B., Information gerrymandering and undemocratic decisions. *Nature*, **573**, 117-121 (2019).
- <sup>91</sup> Ross, L., Greene, D., & House, P., The “false consensus effect”: An egocentric bias in social perception and attribution processes. *Journal of experimental social psychology*, **13**, 279-301 (1977).
- <sup>92</sup> Colleoni, E., Rozza, A., & Arvidsson, A., Echo chamber or public sphere? Predicting political orientation and measuring political homophily in Twitter using big data. *Journal of communication*, **64**, 317-332 (2014).
- <sup>93</sup> Leviston Z., Walker I., Morwinski S., Your opinion on climate change might not be as common as you think. *Nature Climate Change* **3**, 334-337 (2013).
- <sup>94</sup> Baumann, F., Lorenz-Spreen, P., Sokolov, I., Starnini, M., Modeling echo chambers and polarization dynamics in social networks. Accepted for publication in *Physical Review Letters* (2020).
- <sup>95</sup> Sunstein, C. R., The law of group polarization. *Journal of political philosophy*, **10**, 175-195 (2002).
- <sup>96</sup> Sunstein, C. R., *Conspiracy theories and other dangerous ideas*. (Simon and Schuster, 2014).
- <sup>97</sup> Van der Linden, S., The conspiracy-effect: Exposure to conspiracy theories (about global warming) decreases pro-social behavior and science acceptance. *Personality and Individual Differences*, **87**, 171-173 (2015).
- <sup>98</sup> Lewandowsky, S., Oberauer, K., & Gignac, G. E., NASA faked the moon landing—therefore, (climate) science is a hoax: An anatomy of the motivated rejection of science. *Psychological science*, **24**, 622-633 (2013).
- <sup>99</sup> Scheufele, D. A., & Krause, N. M. Science audiences, misinformation, and fake news. *Proc. Natl. Acad. Sci. U.S.A.*, **116**, 7662-7669 (2019).

<sup>100</sup> Lewandowsky, S., Cook, J., Fay, N., & Gignac, G. E. Science by social media: Attitudes towards climate change are mediated by perceived social consensus. *Memory & cognition*, 1-12 (2019).

<sup>101</sup> Muchnik, L., Aral, S., & Taylor, S. J. Social influence bias: A randomized experiment. *Science*, **341**, 647-651 (2013).

<sup>102</sup> [Alipourfard, N., Nettasinghe, B., Abeliuk, A., Krishnamurthy, V., & Lerman, K., Friendship paradox biases perceptions in directed networks. \*Nature Communications\*, \*\*11\*\*, 1-9 \(2020\).](#)

<sup>103</sup> Pennycook G., Rand D. G., Fighting misinformation on social media using crowdsourced judgments of news source quality. *Proc. Natl. Acad. Sci. U.S.A.*, **116**, 2521-2526 (2019).

<sup>104</sup> Ecker, U. K., Lewandowsky, S., & Tang, D. T., Explicit warnings reduce but do not eliminate the continued influence of misinformation. *Memory & cognition*, **38**, 1087-1100 (2010).

<sup>105</sup> Lewandowsky, S., Ecker, U. K., Seifert, C. M., Schwarz, N., & Cook, J. Misinformation and its correction: Continued influence and successful debiasing. *Psychological Science in the Public Interest*, **13**, 106-131 (2012).

<sup>106</sup> <https://about.fb.com/news/2019/10/update-on-election-integrity-efforts/> (accessed: 22.01.2020)

<sup>107</sup> Wineburg, S. & McGrew, S., Lateral reading: Reading less and learning more when evaluating digital information. Working Paper No 2017.A1/Stanford History Education Group, Available at: <https://ssrn.com/abstract=3048994> (2017).

<sup>108</sup> Schultz, P. W., Nolan, J. M., Cialdini, R. B., Goldstein, N. J., & Griskevicius, V., The constructive, destructive, and reconstructive power of social norms. *Psychological science*, **18**, 429-434 (2007).

<sup>109</sup> Hoffrage, U., Lindsey, S., Hertwig, R., & Gigerenzer, G., Communicating statistical information. *Science*, **290**, 2261-2262 (2000).

<sup>110</sup> Tucker, J. A., Theocharis, Y., Roberts, M. E., & Barberá, P., From liberation to turmoil: social media and democracy. *Journal of democracy*, **28**, 46-59 (2017).

<sup>111</sup> <https://www.facebook.com/business/news/insights/capturing-attention-feed-video-creative> (accessed: 08.12.2019)

<sup>112</sup> Kozyreva, A., Lewandowsky, S., & Hertwig, R., Citizens Versus the Internet: Confronting Digital Challenges With Cognitive Tools. Available at: <https://psyarxiv.com/ky4x8/> (2019).

<sup>113</sup> Reijula, S., & Hertwig, R., *Self-nudging and the citizen choice architect*. Accepted for publication in *Behavioural Public Policy* (2020).

---

<sup>114</sup> Noriega-Campero, A., Almaatouq, A., Krafft, P., Alotaibi, A., Moussaid, M., & Pentland, A., The Wisdom of the Network: How Adaptive Networks Promote Collective Intelligence. Available at: <https://arxiv.org/abs/1805.04766> (2018).

<sup>115</sup> Vosoughi, S., Automatic detection and verification of rumors on Twitter. *doctoral dissertation, Massachusetts Institute of Technology*, Cambridge, MA (2015).

<sup>116</sup> Zhou, X., & Zafarani, R., Fake News: A Survey of Research, Detection Methods, and Opportunities. Preprint available at: <https://arxiv.org/abs/1812.00315> (2018).

<sup>117</sup> Martignon, L., Katsikopoulos, K. V., & Woike, J. K., Categorization with limited resources: A family of simple heuristics. *Journal of Mathematical Psychology*, **52**, 352-361 (2008).

<sup>118</sup> Phillips, N. D., Neth, H., Woike, J. K., & Gaissmaier, W., FFTrees: A toolbox to create, visualize, and evaluate fast-and-frugal decision trees. *Judgment and Decision making*, **12**, 344-368 (2017).

<sup>119</sup> Banerjee, S., Chua, A. Y., & Kim, J. J., Don't be deceived: Using linguistic analysis to learn how to discern online review authenticity. *Journal of the Association for Information Science and Technology*, **68**, 1525-1538 (2017).

<sup>120</sup> Cook, J., Lewandowsky, S., Ecker U. K., Neutralizing misinformation through inoculation: Exposing misleading argumentation techniques reduces their influence. *PLOS ONE*, **12**, 1-21 (2017).

<sup>121</sup> Roozenbeek, J., & van der Linden, S., The fake news game: actively inoculating against the risk of misinformation. *Journal of Risk Research*, 1-11 (2018).

<sup>122</sup> Basol, M., Roozenbeek, J., & van der Linden, S., Good News about Bad News: Gamified Inoculation Boosts Confidence and Cognitive Immunity Against Fake News. *Journal of Cognition*, **3**, 2 (2020).

<sup>123</sup> [Dias, N., Pennycook, G., & Rand, D. G., Emphasizing publishers does not effectively reduce susceptibility to misinformation on social media. \*Harvard Kennedy School Misinformation Review\*, \*\*1\*\* \(2020\).](#)